

Title: Text Joins for Data Cleansing and Integration in a Relational Database Management System

Applicants: Koudas et al.

Docket No.: 1209-29

1/12

<pre>INSERT INTO RiIDF(token, idf) SELECT T.token, LOG(S.size)-LOG(COUNT(UNIQUE(*))) FROM RiTokens T, RiSize S GROUP BY T.token, S.size (a) Relation with token idf counts</pre>	<pre>INSERT INTO RiTF(tid, token, tf) SELECT T.tid, T.token, COUNT(*) FROM RiTokens T GROUP BY T.tid, T.token (b) Relation with token tf counts</pre>
<pre>INSERT INTO RiLength(tid, len) SELECT T.tid, SQRT(SUM(I.idf*I.idf*T.tf*T.tf)) FROM RiIDF I, RiTF T WHERE I.token = T.token GROUP BY T.tid (c) Relation with weight-vector lengths</pre>	<pre>INSERT INTO RiWeights(tid, token, weight) SELECT T.tid, T.token, I.idf*T.tf/L.len FROM RiIDF I, RiTF T, RiLength L WHERE I.token = T.token AND T.tid = L.tid (d) Final relation with normalized tuple weight vectors</pre>
<pre>INSERT INTO RiSum(token, total) SELECT R.token, SUM(R.weight) FROM RiWeights R GROUP BY R.token (e) Relation with total token weights</pre>	<pre>INSERT INTO RiSize(size) SELECT COUNT(*) FROM Ri (f) Dummy relation used to create RiIDF</pre>

Fig. 1

• Title: Text Joins for Data Cleansing and Integration in a Relational Database Management System
Applicants: Koudas et al.
Docket No.: 1209-29

2/12

```
SELECT r1w.tid AS tid1, r2w.tid AS tid2
FROM R1Weights r1w, R2Weights r2w
WHERE r1w.token = r2w.token
GROUP BY r1w.tid, r2w.tid
HAVING SUM(r1w.weight*r2w.weight) ≥ φ
```

Fig. 2

Title: Text Joins for Data Cleansing and Integration in a Relational Database Management System

Applicants: Koudas et al.

Docket No.: 1209-29

3/12

```
SELECT rw.tid, rw.token, rw.weight/rs.total AS P
FROM   RiWeights rw, RiSum rs
WHERE  rw.token = rs.token
```

Fig 3

Title: Text Joins for Data Cleansing and Integration in a Relational Database Management System

Applicants: Koudas et al.

Docket No.: 1209-29

4/12

```
INSERT INTO RiSample(tid,token,c)
SELECT rw.tid, rw.token, ROUND(S * rw.weight/rs.total, 0) AS c
FROM   RiWeights rw, RiSum rs
WHERE  rw.token = rs.token
```

Fig. 4

Title: Text Joins for Data Cleansing and Integration in a Relational Database Management System

Applicants: Koudas et al.

Docket No.: 1209-29

5/12

```
SELECT r1w.tid AS tid1, r2s.tid AS tid2
FROM Riweights r1w, R2sample r2s, R2sum r2sum, R1V r1v
WHERE r1w.token = r2s.token AND r1w.token = r2sum.token AND r1w.tid = r1v.tid
GROUP BY r1w.tid, r2s.tid, r1v.Tv
HAVING SUM(r1w.weight * r2sum.total / r1v.Tv) ≥ S * φ' / r1v.Tv
```

Fig. 5

Title: Text Joins for Data Cleansing and Integration in a Relational Database Management System

Applicants: Koudas et al.

Docket No.: 1209-29

6/12

```
SELECT tid1, tid2
FROM
(
  SELECT r1w.tid AS tid1, r2s.tid AS tid2, SUM(r1w.weight * r2sum.total) AS Ci
  FROM Riweights r1w, R2sample r2s, R2sum r2sum
  WHERE r1w.token = r2s.token AND r1w.token = r2sum.token AND r1w.tid = r1v.tid
  GROUP BY r1w.tid, r2s.tid
  UNION ALL
  SELECT r1s.tid AS tid1, r2w.tid AS tid2, SUM(r2w.weight * risum.total) AS Ci
  FROM R2weights r2w, Risample r1s, Risum risum
  WHERE r2w.token = r1s.token AND r2w.token = risum.token AND r2w.tid = r2v.tid
  GROUP BY r2w.tid, r1s.tid
) SYM
GROUP BY tid1, tid2
HAVING AVG(Ci)  $\geq$  S *  $\phi'$ 
```

Fig. 6

Title: Text Joins for Data Cleansing and Integration in a Relational Database Management System

Applicants: Koudas et al.

Docket No.: 1209-29

7/12

```
SELECT  ris.tid AS tid1, r2s.tid AS tid2
FROM    R1Sample ris, R2Sample r2s, R1Sum risum, R2Sum r2sum
WHERE   ris.token = risum.token AND R2Sample.token = r2sum.token AND ris.token = r2s.token
GROUP BY ris.tid, r2s.tid
HAVING  SUM(risum.total * r2sum.total) ≥ S * S * φ'
```

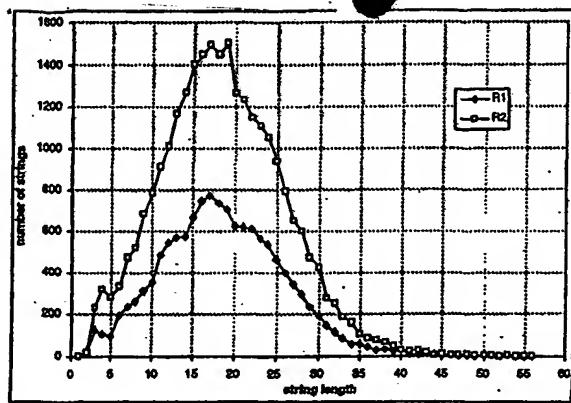
Fig. 7

Title: Text Joins for Data Cleansing and Integration in a Relational Database Management System

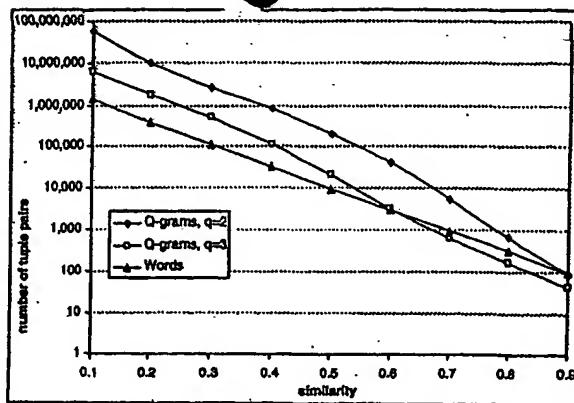
Applicants: Koudas et al.

Docket No.: 1209-29

8/12

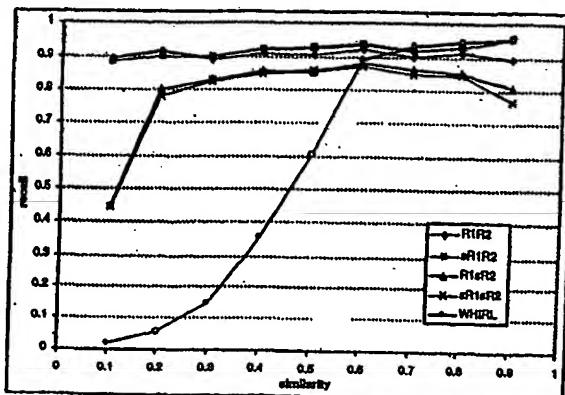
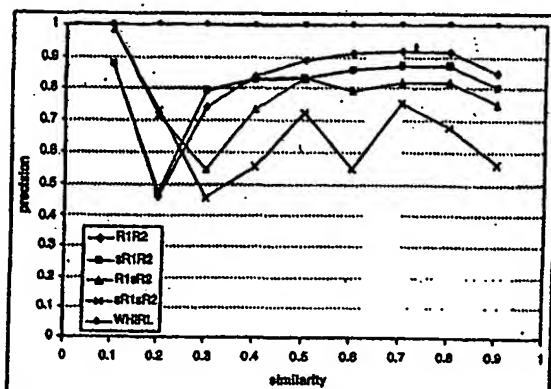


(a) String lengths in data sets R_1 and R_2 .

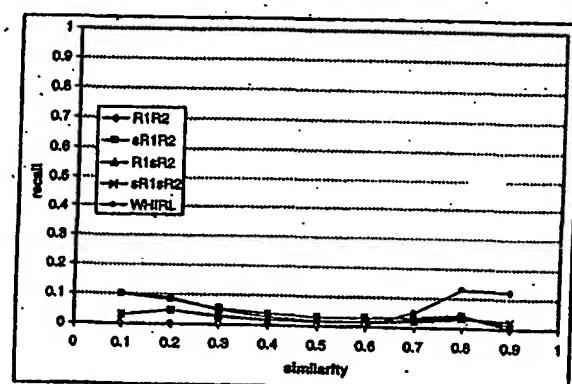
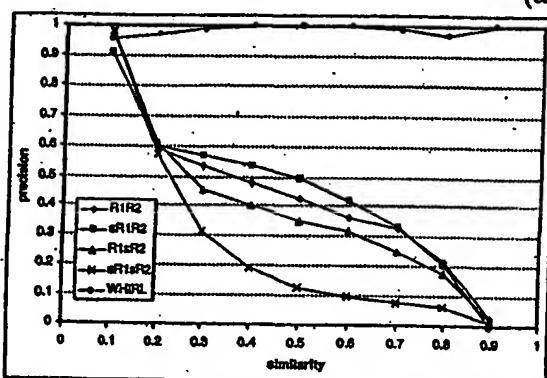


(b) The size of $R_1 \bowtie_{\phi} R_2$ for different similarity thresholds and token choices.

Fig 8



(a) Words



(b) Q -grams with $q = 2$

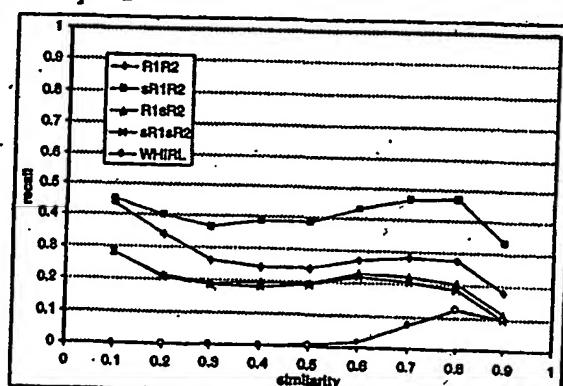
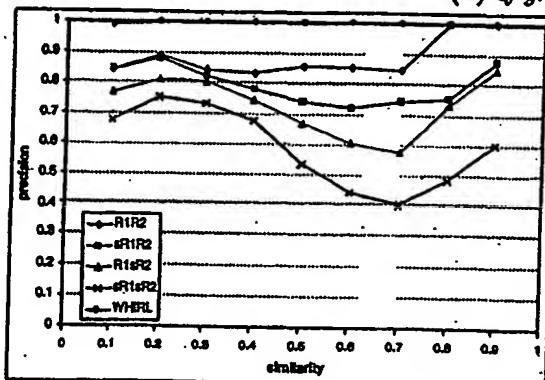


Fig. 9

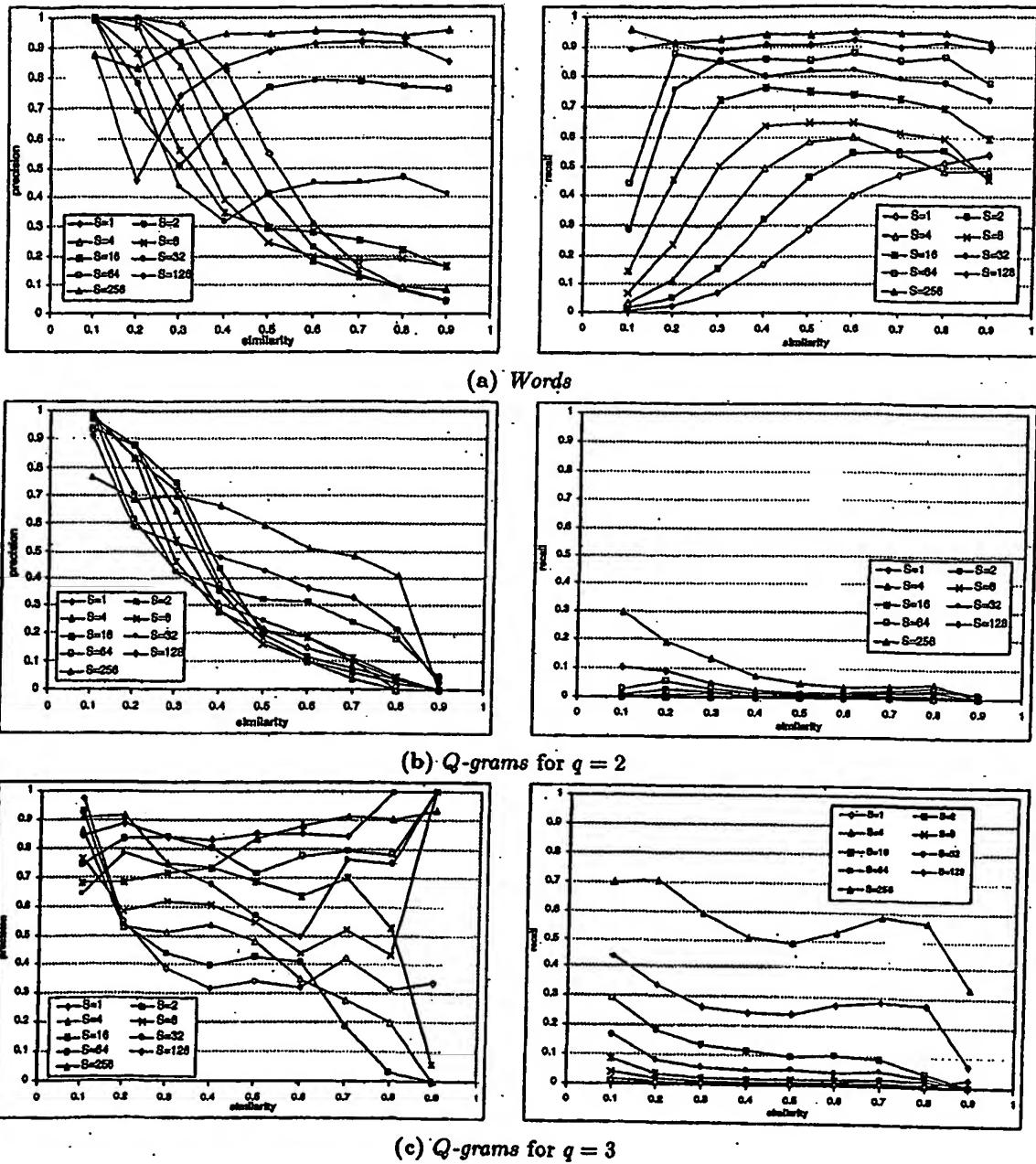


Fig. 10

Title: Text Joins for Data Cleansing and Integration in a Relational Database Management System

Applicants: Koudas et al.

Docket No.: 1209-29

11/12

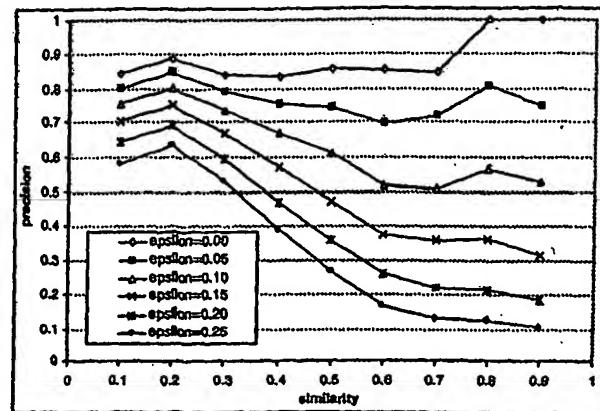
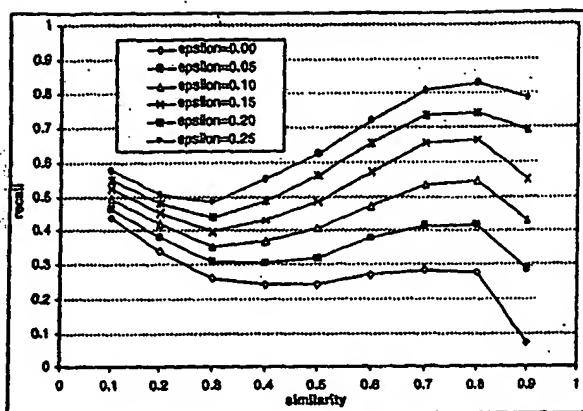
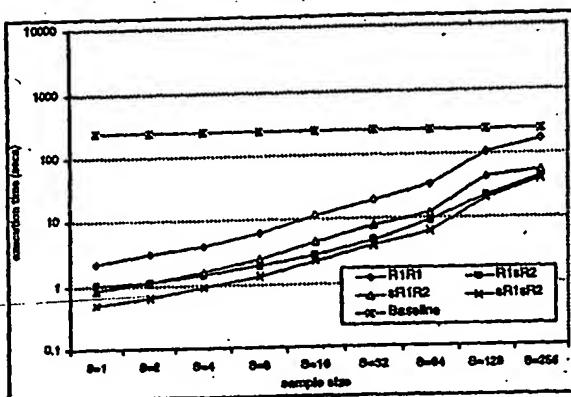
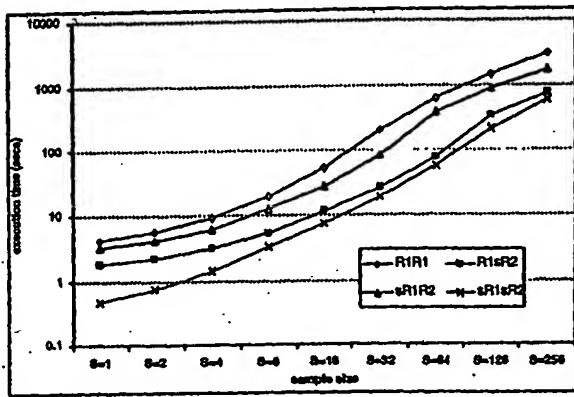
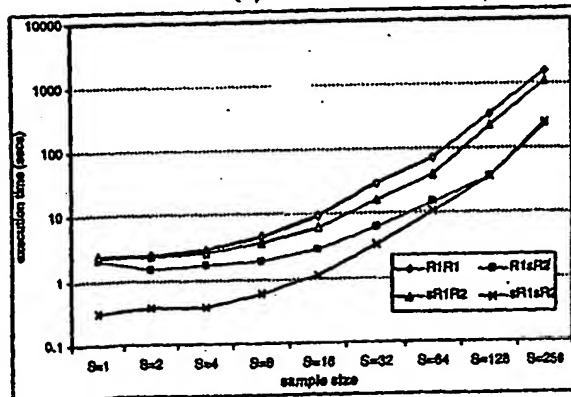
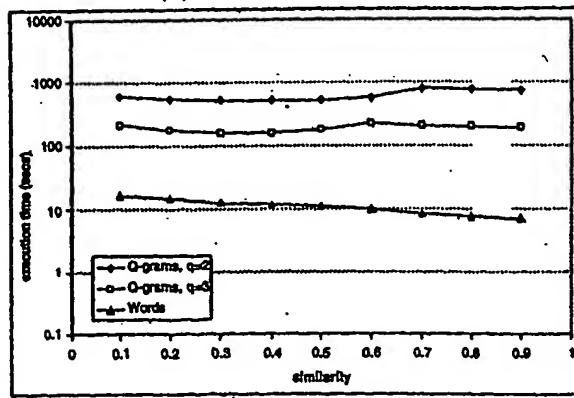


Fig. 11



(a) Words

(b) Q -grams with $q = 3$ (c) Q -grams with $q = 2$ 

(d) WHIRL

Fig. 12